

NumCSE exercise sheet 4

Convolution and FFT, filtering, Lagrange interpolation

alexander.dabrowski@sam.math.ethz.ch
oliver.rietmann@sam.math.ethz.ch

November 21, 2018

Exercise 4.1. (Long exercise) Sparse vectors, their convolution and FFT

Given a vector \mathbf{x} of numbers, we want to represent it with a new data structure which should be suitable in case most of its entries are negligible. We also want to implement convolution and fast Fourier transform for this new data structure.

(a) Consider the following two **structs**:

- **template<class T> duplet**, which has members
 - **int ind**, the index of an element in \mathbf{x} ;
 - **T val**, the value of the element at position **ind**;
- **template<class T> sparse_vec**, which has members
 - **double tol**. We consider negligible any number with absolute value smaller than **tol**;
 - **duplets**, a **vector** of objects **duplet<T>** which stores the indices and values of the non-negligible elements of \mathbf{x} ;
 - **len**, an integer which records the length of \mathbf{x} .

Implement in **sparse_vec<T>** the following utility methods:

- **void append(int ind, T val)**, which appends to **duplets** a new duplet with index **ind** and value **val**, if $\text{abs}(\text{val}) \geq \text{tol}$;
- **void cleanup()**, which sorts **duplets** with respect to **ind**, eliminates repetitions, with the convention that if two elements have the same index their values should be added up, and eliminates any element with norm smaller than **tol** or with index larger than **len-1**. If $\mathbf{x}.\text{cleanup}()$ leaves \mathbf{x} as it is, we say that \mathbf{x} is **clean**;
- **T get_val(int ind)**, which returns the element of \mathbf{x} in position **ind**. It should work in $O(\log \text{duplets.size}())$ time and assumes that \mathbf{x} is clean.

Solution:

```
13 template<class T> struct duplet {  
14     int ind;  
15     T val;  
16  
17     duplet(int p, T v) {  
18         ind = p;  
19         val = v;  
20     }  
21};
```

```

22
23 template<class T> struct sparse_vec {
24     double tol = 1e-6;
25     vector<duplet<T>> duplets;
26     int len=0;
27
28     sparse_vec(int l) {
29         len = l;
30     }
31
32     void append(int ind, T val) {
33         if (abs(val) >= tol) {
34             duplet<T> newduplet(ind, val);
35             duplets.push_back(newduplet);
36         }
37     }
38
39     void cleanup() {
40         std::sort(duplets.begin(), duplets.end(),
41                 [] (duplet<T> x, duplet<T> y) { return x.ind < y.ind; });
42
43         vector<duplet<T>> newduplets;
44         T tmp = 0;
45         for (int i=0; i<duplets.size(); i++) {
46             if (i == duplets.size()-1 || duplets[i+1].ind != duplets[i].ind) {
47                 tmp += duplets[i].val;
48                 if (abs(tmp) >= tol) {
49                     duplet<T> newduplet(duplets[i].ind, tmp);
50                     newduplets.push_back(newduplet);
51                 }
52                 tmp = 0;
53             }
54             else {
55                 tmp += duplets[i].val;
56             }
57         }
58         duplets = newduplets;
59     }
60
61     T get_val(int ind) const {
62         if (duplets.empty())
63             return 0;
64         return _get_val(ind, 0, duplets.size());
65     }
66
67     T _get_val(int ind, int n1, int n2) const {
68         if (n2 < n1)
69             return 0;
70         int m = (n1 + n2)/2;
71         if (duplets[m].ind == ind)
72             return duplets[m].val;
73         if (duplets[m].ind < ind)
74             return _get_val(ind, m+1, n2);
75         if (duplets[m].ind > ind)
76             return _get_val(ind, n1, m-1);

```

Example: A representation of $x = (0, 0, 1, 0, 7 + i, 0)$ can be obtained with a `sparse_vec` with `duplets = {(ind=4, val=6-I), (ind=1, val=1e-7), (ind=22, val=9), (ind=4, val=1+2I), (ind=2, val=1)}`, `len=6`, `tol=1e-6`. Notice that `duplets.size() != len`. Running `cleanup()` on such an object, `duplets` would become `{(ind=2, val=1), (ind=4, val=7+I)}`.

Important notes: In all the subproblems below, you can assume all vectors of type `sparse_vec` which are passed as inputs are already clean. Do not convert the `sparse_vecs` to/from dense vectors.

- (b) In struct `sparse_vec` implement a method `cwise_mult` which returns the component-wise multiplication between `sparse_vec a` and `sparse_vec b` in $O(sz_a + sz_b)$ time, where sz_a , sz_b are respectively the number of non-negligible elements of `a` and `b`.

Solution:

```

79     static sparse_vec cwise_mult(const sparse_vec &a, const sparse_vec &b) {
80         sparse_vec out(max(a.len,b.len));
81         int ia = 0, ib = 0;
82         while(ia < a.duplets.size() && ib < b.duplets.size()) {
83             int inda = a.duplets[ia].ind, indb = b.duplets[ib].ind;
84             if (inda >= a.len || indb >= b.len) break;
85             if (inda == indb) {
86                 out.append(inda, a.duplets[ia].val * b.duplets[ib].val);
87                 ia++;
88                 ib++;
89             }
90             else if (inda < indb)
91                 ia++;
92             else if (inda > indb)
93                 ib++;
94         }
95         return out;
96     }

```

- (c) Implement a method `conv` which returns the discrete convolution of `sparse_vec a` and `sparse_vec b` in $O(sz_a sz_b)$ time.

Solution:

```

98     static sparse_vec conv(const sparse_vec &a, const sparse_vec &b) {
99         sparse_vec out(a.len + b.len - 1);
100        for (auto x : a.duplets) {
101            for (auto y : b.duplets) {
102                out.append(x.ind + y.ind, x.val * y.val);
103            }
104        }
105        return out;
106    }

```

- (d) Implement a method `fft` which returns the discrete Fourier transform of `sparse_vec x` in $O(n(\log n)^2)$ time, where $n = x.len$. You can assume that n is a power of 2. Your function should return a `sparse_vec` which is clean.

(Optional) Improve your code so that it runs in $O(n \log n)$. Even if you don't implement this improvement, you can suppose in the following subproblems that the runtime of `fft` is $O(n \log n)$.

Solution:

```

108 static sparse_vec fft(const sparse_vec &x) {
109     int n = x.len;
110     if (n<=1) return x;
111
112     sparse_vec even(n/2);
113     sparse_vec odd(n/2);
114     for (auto duplet : x.duplets) {
115         if (duplet.ind % 2 == 0)
116             even.append(duplet.ind/2, duplet.val);
117         else
118             odd.append((duplet.ind-1)/2, duplet.val);
119     }
120
121     sparse_vec f0 = fft(even);
122     sparse_vec f1 = fft(odd);
123
124     T omega = exp(-2.*PI/n*I);
125     T s(1.,0.);
126     sparse_vec tot(n);
127
128     /* O(n*(log(n))^2) solution
129     for (int k=0; k<n; k++) {
130         tot.append(k, f0.get_val(k % (n/2)) + f1.get_val(k % (n/2))*s);
131         s *= omega;
132     }
133     */
134
135     /* Another O(n*(log(n))^2) solution
136     for (auto p : f0.duplets) {
137         tot.append(p.ind, p.val);
138         tot.append(p.ind + n/2, p.val);
139     }
140     for (auto p : f1.duplets) {
141         tot.append(p.ind, p.val * std::pow(omega, p.ind));
142         tot.append(p.ind + n/2, p.val * std::pow(omega, p.ind + n/2));
143     }
144
145     tot.cleanup();
146     */
147
148     auto merge_insert = [&] (int shift) {
149         int i0 = 0, i1 = 0;
150         int f0_size = f0.duplets.size(), f1_size = f1.duplets.size();
151         while (i0 < f0_size && i1 < f1_size) {
152             int ind0 = f0.duplets[i0].ind + shift;
153             int ind1 = f1.duplets[i1].ind + shift;
154             if (ind0 == ind1) {
155                 T tmp = 0;
156                 do {
157                     tmp += f0.duplets[i0].val;
158                     tmp += f1.duplets[i1].val * pow(omega, ind0);
159                     i0++; i1++;

```

```

160     } while (i0 < f0_size && i1 < f1_size &&
161         f0.duplets[i0].ind == ind0 && f1.duplets[i1].ind == ind0);
162     tot.append(ind0, tmp);
163 }
164 else if (ind0 < ind1) {
165     tot.append(ind0, f0.duplets[i0].val);
166     i0++;
167 }
168 else if (ind0 > ind1) {
169     tot.append(ind1, f1.duplets[i1].val * pow(omega, ind1));
170     i1++;
171 }
172 }
173 while(i0 < f0_size) {
174     int ind0 = f0.duplets[i0].ind + shift;
175     tot.append(ind0, f0.duplets[i0].val);
176     i0++;
177 }
178 while(i1 < f1_size) {
179     int ind1 = f1.duplets[i1].ind + shift;
180     tot.append(ind1, f1.duplets[i1].val * pow(omega, ind1));
181     i1++;
182 }
183 };
184 merge_insert(0);
185 merge_insert(n/2);
186
187 return tot;
188 }
```

sparse_conv_fft.cpp

- (e) Implement a method `ifft` which returns the inverse discrete Fourier transform of `sparse_vec` `x` in $O(n \log n)$ time. You can assume that n is a power of 2.

Solution:

```

190 static sparse_vec ifft(const sparse_vec &x) {
191     double n = x.len;
192     sparse_vec out(n);
193     sparse_vec x_conj(n);
194     for (auto duplet : x.duplets) {
195         x_conj.append(duplet.ind, std::conj(duplet.val));
196     }
197     for (auto duplet : fft(x_conj).duplets) {
198         out.append(duplet.ind, std::conj(duplet.val)/n);
199     }
200     return out;
201 }
```

sparse_conv_fft.cpp

- (f) Let `a` and `b` be two complex vectors of length $n + 1$ and n respectively. Assume that n is a power of 2. Implement a method `conv_fft` which returns the discrete convolution of `a` and `b` in $O(n \log n)$ time.

Solution:

```

203     static sparse_vec conv_fft(sparse_vec a, sparse_vec b) {
204         int N = a.len + b.len - 1;
205         a.len = N;
206         b.len = N;
207         return ifft(cwise_mult(fft(a), fft(b)));
208     }

```

sparse_conv_fft.cpp

- (g) Can `conv_fft(x, y)` be asymptotically slower than `conv(x, y)` for a particular choice of `x` and `y`? Briefly motivate your answer.

Solution: Yes, for instance if `x` and `y` are very sparse but their discrete Fourier transforms are not. For example if `x` and `y` have all zero elements except for a single non-negligible 1 in the first position, their `fft` are the constant 1 vectors. Then `conv_fft` must at least calculate all these ones, while `conv` will just multiply together two numbers to get the result.

Exercise 4.2. 2D convolution, FFT2 and Laplace filter

We review the two dimensional convolution, its relation with the discrete Fourier transform, and implement a discrete Laplacian filter.

Consider the 2-dimensional infinite arrays $X = (X_{k_1, k_2})_{k_1, k_2 \in \mathbb{Z}}$ and $Y = (Y_{k_1, k_2})_{k_1, k_2 \in \mathbb{Z}}$. We can define their convolution as the 2-dimensional infinite array $X * Y$ such that

$$(X * Y)_{k_1, k_2} = \sum_{j_1=-\infty}^{\infty} \sum_{j_2=-\infty}^{\infty} X_{j_1, j_2} Y_{k_1 - j_1, k_2 - j_2}.$$

In the same way as in the 1-dimensional case, given two matrices $X \in \mathbb{R}^{m_1, m_2}$ and $Y \in \mathbb{R}^{n_1, n_2}$ we can define their *discrete convolution* as the convolution between the zero extensions of X and Y trimmed of unnecessary zeros, and their *circular convolution* as the smallest period of the convolution between the periodic extension of X and the zero extension of Y .

Consider two matrices $A \in \mathbb{R}^{n, m}$ and $F \in \mathbb{R}^{k, k}$ with $k < n, m$.

- (a) How should we extend A and F to larger matrices \tilde{A} and \tilde{F} so that the discrete convolution of A with F is equal to the circular convolution of \tilde{A} with \tilde{F} ?

Solution: By padding A and F with zeros so that the dimensions of \tilde{A} and \tilde{F} are $L_r \times L_c$ where $L_r = (\text{number of rows of } A) + (\text{number of rows of } F) - 1$ and $L_c = (\text{number of columns of } A) + (\text{number of columns of } F) - 1$.

- (b) By recalling to the 1-dimensional fast fourier transform (see the Eigen::FFT module), implement a function `fft2` which returns the 2-dimensional fast fourier transform of a matrix.

Solution:

```

9 template <typename Scalar> void fft2(MatrixXcd &C, const MatrixBase<Scalar> &Y) {
10    int m = Y.rows(), n = Y.cols();
11    C.resize(m,n);
12    MatrixXcd tmp(m,n);
13
14    FFT<double> fft; // Helper class for DFT
15    // Transform rows of matrix Y
16    for (int k = 0; k < m; k++) {
17        VectorXcd tv(Y.row(k));
18        tmp.row(k) = fft.fwd(tv).transpose();
19    }
20
21    // Transform columns of temporary matrix
22    for (int k = 0; k < n; k++) {
23        VectorXcd tv(tmp.col(k));
24        C.col(k) = fft.fwd(tv);
25    }
26}

```

fft2.cpp

- (c) Implement a function `ifft2` which returns the 2-dimensional inverse fast fourier transform of a matrix.

Solution:

```

28 template <typename Scalar> void ifft2(MatrixXcd &C, const MatrixBase<Scalar> &Y) {
29    int m = Y.rows(), n = Y.cols();
30    fft2(C, Y.conjugate());
31    C = C.conjugate()/(m*n);

```

- (d) Implement an efficient function with arguments A and F which computes their discrete convolution.

Hint: use the result derived in the previous steps and the two dimensional circular convolution theorem.

Solution:

```

34 void conv2(MatrixXcd &C, const MatrixXcd &A1, const MatrixXcd &A2) {
35     // returns discrete convolution between A1 and A2 (fft implementation)
36     int n1 = A1.rows();
37     int m1 = A1.cols();
38     int n2 = A2.rows();
39     int m2 = A2.cols();
40     int Lrow = n1+n2-1;
41     int Lcol = m1+m2-1;
42
43     MatrixXcd A1_ext = MatrixXcd::Zero(Lrow, Lcol);
44     MatrixXcd A2_ext = MatrixXcd::Zero(Lrow, Lcol);
45     MatrixXcd tmp1 = MatrixXcd::Zero(Lrow, Lcol);
46     MatrixXcd tmp2 = MatrixXcd::Zero(Lrow, Lcol);
47     A1_ext.topLeftCorner(n1,m1) = A1;
48     A2_ext.topLeftCorner(n2,m2) = A2;
49
50     fft2(tmp1, A1_ext);
51     fft2(tmp2, A2_ext);
52     ifft2(C, tmp1.cwiseProduct(tmp2));
53 }
```

- (e) The discrete Laplacian filter is defined as

$$F = \begin{pmatrix} 0 & -1 & 0 \\ -1 & 4 & -1 \\ 0 & -1 & 0 \end{pmatrix}.$$

It is often used to detect edges in images. Test your filter on the black and white “picture” provided in the template.

Exercise 4.3. Lagrange interpolation.

Fix $n \in \mathbb{N}_0$ and let $t_0, \dots, t_n \in \mathbb{R}$ be distinct nodes, i.e. $t_i \neq t_j$ if $i \neq j$. Then

$$L_i(t) := \prod_{\substack{j=0 \\ j \neq i}}^n \frac{t - t_j}{t_i - t_j}$$

is called the i -th *Lagrange polynomial* for these given nodes. For $y_0, \dots, y_n \in \mathbb{R}$ we call

$$p(t) := \sum_{i=0}^n y_i L_i(t)$$

the *Lagrange interpolant* through $(t_i, y_i)_{i=0}^n$.

(a) Prove that

$$\sum_{i=0}^n L_i(t) = 1$$

for all $t \in \mathbb{R}$.

Hint: Choose $y_i = 1$ for all $i \in \{0, \dots, n\}$ and use uniqueness of Lagrange interpolants.

Solution: Consider the interpolation points $(t_i, 1)_{i=0}^n$. Clearly, an interpolating polynomial for these points is $p \equiv 1$. By uniqueness of Lagrange interpolants, we have

$$1 = p(t) = \sum_{i=0}^n L_i(t)$$

for all $t \in \mathbb{R}$.

(b) For $t \in \mathbb{R}$ define $\omega(t) := \prod_{j=0}^n (t - t_j)$ and fix $i \in \{0, \dots, n\}$. Prove that $\omega'(t_i) \neq 0$ and

$$L_i(t) = \omega(t) \frac{\lambda_i}{t - t_i}$$

for all $t \in \mathbb{R}$, where

$$\lambda_i := \frac{1}{\omega'(t_i)}.$$

Solution: Fix $i \in \{0, \dots, n\}$. By the product rule, we have

$$\omega'(t) = \frac{d}{dt}(t - t_i) \prod_{\substack{j=0 \\ j \neq i}}^n (t - t_j) = \prod_{\substack{j=0 \\ j \neq i}}^n (t - t_j) + (t - t_i) \frac{d}{dt} \prod_{\substack{j=0 \\ j \neq i}}^n (t - t_j)$$

for all $t \in \mathbb{R}$. In particular for $t = t_i$, we obtain

$$\omega'(t_i) = \prod_{\substack{j=0 \\ j \neq i}}^n (t_i - t_j) \neq 0$$

since the nodes are distinct. Therefore,

$$L_i(t) = \frac{\prod_{j=0}^n (t - t_j)}{\prod_{\substack{j=0 \\ j \neq i}}^n (t_i - t_j)} = \frac{\prod_{j=0}^n (t - t_j)}{(t - t_i) \omega'(t_i)} = \omega(t) \frac{\lambda_i}{t - t_i}$$

for all $t \in \mathbb{R}$.

i	t_i	y_i
0	-1	2
1	0	-4
2	1	6

Table 1: data

- (c) Compute the Lagrange interpolant $p(t)$ corresponding to the data given in Table 1.

Solution: The Lagrange polynomials are given by

$$\begin{aligned}L_0(t) &= \frac{t^2 - t}{2} \\L_1(t) &= 1 - t^2 \\L_2(t) &= \frac{t^2 + t}{2}\end{aligned}$$

yielding the interpolant

$$p(t) = t^2 - t + 4(t^2 - 1) + 3(t^2 + t) = 8t^2 + 2t - 4.$$

- (d) Use the Newton basis approach to compute the interpolating polynomial $\tilde{p}(t)$ for the data in Table 1.

Solution: We have

$$\tilde{p}(t) = \sum_{i=0}^2 a_i N_i(t)$$

for all $t \in \mathbb{R}$, where

$$\begin{aligned}N_0(t) &= 1 \\N_1(t) &= (t - t_0) \\N_2(t) &= (t - t_0)(t - t_1)\end{aligned}$$

and

$$\begin{aligned}a_0 &= 2 \\a_1 &= -6 \\a_2 &= 8.\end{aligned}$$

Hence we obtain

$$\tilde{p}(t) = 2 - 6(t + 1) + 8t(t + 1) = 8t^2 + 2t - 4.$$

- (e) Is $\tilde{p}(t)$ different from $p(t)$? Explain your answer.

Solution: We have $\tilde{p}(t) = p(t)$ for all $t \in \mathbb{R}$ because the interpolating polynomial is unique.

- (f) What are the advantages of using the Newton basis compared to the Lagrange polynomials?

Solution: The Lagrange interpolant might become numerically unstable if some of the nodes t_i are close to each other, as it involves computing $t_i - t_j$ for $i \neq j$. In the Newton basis, we only compute $t - t_j$, so it is numerically more stable. Moreover, inserting a new node does not affect the previous coefficients in the Newton basis (while it does in the Lagrange basis).